

**ETH login ID:**

(Please print in capital letters)

--	--	--	--	--	--	--	--	--	--	--	--

**Full name:**

---

**263-0007-00L: Advanced Systems Lab**

ETH Computer Science, Spring 2020

Midterm Exam

Monday, July 08, 2020

**Instructions**

- Write your full name and login ID on the front.
- Make sure that your exam is not missing any sheets.
- No extra sheets are allowed.
- The exam has a maximum score of 100 points.
- No books, notes, calculators, laptops, cell phones, or other electronic devices are allowed.

Problem 1 ( $16 = 2+4+2+4+4$ )	<input type="text"/>
Problem 2 ( $15 = 6+6+3$ )	<input type="text"/>
Problem 3 ( $18 = 2+6+6+4$ )	<input type="text"/>
Problem 4 ( $14 = 2+2+3+5+2$ )	<input type="text"/>
Problem 5 ( $18 = 2+2+2+4+2+6$ )	<input type="text"/>
Problem 6 ( $19 = 2+2+2+2+4+2+5$ )	<input type="text"/>
<hr/>	
<b>Total (100)</b>	<input type="text"/>

## Problem 1: Operational Intensity (16 = 2+4+2+4+4)

Consider the computation  $y^\top = (Ax)^\top A$  where  $x, y$  are column vectors of doubles of length  $n$  and  $A$  is an  $n \times n$  matrix of doubles stored in row major order. No temporary arrays are used in this computation. `sizeof(double) = 8`. Make the following assumptions:

- A write-back/write-allocate and initially empty (cold) cache of size  $\gamma$  bytes.
- $n$  is a multiple of the cache block of size  $B = 64$  bytes.
- Arrays are aligned with the cache blocks.

**In the derivations you can omit lower order terms** (writing  $\approx$  instead of  $=$ ). Show your work.

1. Determine the flop count  $W(n)$  of the computation.
2. Determine an upper bound for the operational intensity  $I(n)$  of the computation considering only compulsory data movement. Consider only reads and ignore writes.
3. For which sizes  $n$  would you expect the bound for  $I(n)$  calculated in 2) to be roughly accurate (i.e., roughly equal to  $I(n)$ )?

4. Consider now that the computation is performed on a submatrix  $A'$  and subvector  $x'$  obtained by accessing  $A$  and  $x$  with a stride  $k$ . Assume that  $n$  is divisible by  $k$ . In Matlab notation, this is equivalent to  $y'^T = (A'x')^T A'$  where  $A' = A_{0:k:n^2-1}$ ,  $x' = x_{0:k:n-1}$ , and  $y' = y_{0:k:n-1}$ , so  $A'$  is an  $n \times \frac{n}{k}$  matrix. For convenience, we also include the code of the computation:

```
1 /* Arrays y and tmp are initialized with zeros */
2 void f(double *A, double * x, double *y, double *tmp, int n, int k) {
3
4     // tmp = A*x (with stride k)
5     for (int i = 0; i < n; ++i)
6         for (int j = 0; j < n; j += k)
7             tmp[i] += A[i*n + j]*x[j];
8
9     // y = tmp^T*A (with stride k)
10    for (int j = 0; j < n; j += k)
11        for (int i = 0; i < n; ++i)
12            y[j] += tmp[i]*A[i*n + j];
13 }
```

Determine a new upper bound for the operational intensity  $I(n)$  for the following cases considering only compulsory data movement. Consider only reads and ignore writes. Briefly explain your answer.

(a) For  $k = 7$

(b) For  $k = 16$



- Using the bound calculated in 1) and assuming  $k = 10n$ , determine the values of  $\beta$  for which the Sparse MVM using CSR format is guaranteed to be memory bound. Show your derivation.

### Problem 3: Cache Mechanics (18 = 2+6+6+4)

Consider a 2-way set associative cache with a block size of 16 bytes and 4 sets. The replacement policy is LRU. The write policy is write-back/write-allocate. The cache is initially cold. Consider the following code:

```
1 void compute(double *A) {
2     double t1;
3     for (int i = 0; i < 8; ++i) {
4         t1 = A[i*s1];
5         A[i*s2] = t1 + 2.0;
6     }
7 }
```

Assume that A starts at memory address 0. All memory accesses happen in exactly the order that they appear and are inbound. Variables `t1` and `i` are kept in registers. Answer the following. Show your work. Hint: It helps to draw the cache.

1. What is the size of the cache in bytes?
2. For each of the following values of  $s1$  and  $s2$  and considering read and write misses do the following three things: i) determine the hit/miss pattern (something like HMMHM...); ii) draw the state of the cache after the loop has executed; iii) state what kind(s) of locality does this execution have.
  - (a)  $s1 = 2$  and  $s2 = 4$ ?

(b)  $s_1 = 3$  and  $s_2 = 5$ ?

3. What is the maximum number of cache misses for this code? Find a combination of  $s_1$  and  $s_2$  that achieves this maximum with  $2 \leq s_1 \leq s_2 \leq 9$  and where all misses are compulsory misses.

## Problem 4: ILP ( $14 = 2+2+3+5+2$ )

Consider the AVX code below:

```
1 #include <immintrin.h>
2
3 /* Assume a, b and c are of size 4*n */
4 void f (double *a, double *b, double *c, unsigned int n) {
5     __m256d va, vb, v1, v2, v3, v4, v5;
6     size_t i;
7
8     for (i = 0; i < n; i++) {
9         va = _mm256_loadu_pd    (a + 4*i);
10        vb = _mm256_loadu_pd    (b + 4*i);
11
12        v1 = _mm256_add_pd      (va, vb);
13        v2 = _mm256_add_pd      (v1, vb);
14
15        v3 = _mm256_mul_pd      (va, va);
16        v4 = _mm256_permute_pd  (v3, 0x5);
17
18        v5 = _mm256_add_pd      (v2, v4);
19
20        _mm256_storeu_pd       (c + 4*i, v5);
21    }
22 }
```

Assume that the working set fits in L1 cache and is already loaded (warm cache). Further, assume that the code is executed on a machine with the following latency, throughput, and port information:

Instruction	Latency [cycles]	Throughput (all ports) [instructions/cycle]	Port
<code>_mm256_storeu_pd</code>	1	1	4
<code>_mm256_loadu_pd</code>	1	2	2/3
<code>_mm256_permute_pd</code>	2	1	5
<code>_mm256_mul_pd</code>	5	2	0/1
<code>_mm256_add_pd</code>	4	1	1

Show your derivations or briefly explain your solutions.

1. Determine the flop count  $W(n)$  of  $f$ .





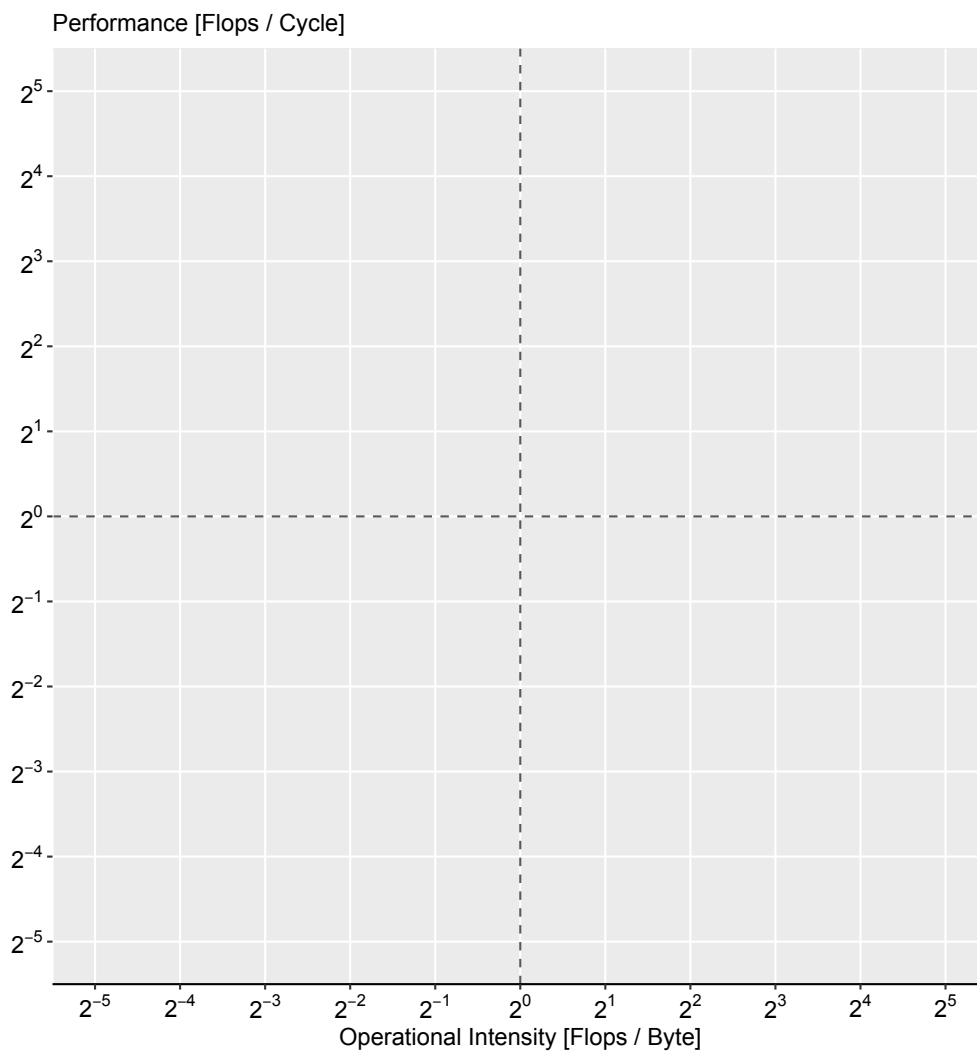
4. Assume that we unroll the loop by a factor  $K$  and reorder the instructions such that all the  $2K$  loads occur at the beginning of the new (unrolled) loop body, followed by the computations ( $4K$  arithmetic SIMD instructions and  $K$  permutes), and finally the  $K$  stores. Assume that register pressure is not an issue (i.e., all loop variables can be held in registers). Take latency, throughput and dependency information into account and answer the following:
- (a) Estimate the runtime (latency) as a function of  $K$  of one loop iteration of this new unrolled loop.

- (b) Based on the latency of one iteration determined in 4a), estimate a lower bound (as tight as possible) for the runtime of the unrolled version of  $\text{f}$ . Give the bound as a function of  $K$  and  $n$ . Assume that  $n$  is divisible by  $K$  and that loads and stores are always computed “in order” by the processor, meaning that a load from one iteration cannot be scheduled before a store from a previous iteration.

## Problem 5: Roofline (18 = 2+2+2+4+2+6)

Assume a computer with the following features:

- A CPU that supports single precision floating point operations. The relevant port information is as follows:  
Port 1: fma, mul, add.  
Port 2: fma, mul.
- Each of these operations has a throughput of 1 per port.
- It does not support any SIMD operations.
- A write-back/write-allocate cache of size 1 MB with cache block size  $B = 64$  bytes.
- The read (memory) bandwidth is 4 floats per cycle. `sizeof(float) = 4`.



1. Draw the roofline plot for this computer into the above graph. Annotate the lines so we see your reasoning.
2. Consider the following code which computes the sum of three vectors  $x$ ,  $y$ , and  $z$  (each has size  $n$ ). Assume that  $x$ ,  $y$ , and  $z$  are allocated sequentially in memory, one after the other, and the first element of  $x$  maps to the start of the first cache block:

```
1 void compute(float *x, float *y, float *z, int n) {  
2     for(int i=0; i < n; i++){  
3         x[i] = x[i] + y[i] + z[i];  
4     }  
5 }
```

Based on the instruction mix (i.e. ignoring dependencies), which performance is maximally achievable for this function and why? Draw an associated tighter horizontal roofline into the plot above. Assume that no FMA instructions are used.

3. At what operational intensity  $I(n)$  does this new horizontal roofline intersect with the memory roofline?



## Problem 6: Sampler (19 = 2+2+2+2+4+2+5)

1. Assume that the peak performance  $\pi$  of computers doubles once every 4 years, whereas the memory bandwidth  $\beta$  doubles once every 8 years. Does this mean that more computations tend to become compute bound or memory bound in the sense defined by the roofline plot? Justify your answer.
2. You are given a double type square matrix  $A \in \mathbb{R}^{n \times n}$  with  $k$  non-zero entries. What is the largest  $k$  that makes it possible to save storage space by representing it in compressed sparse row (CSR) format, as opposed to the standard dense representation? Assume that indices are represented by (4-byte) integers. Show your work.
3. What is the advantage of using a write-back/write-allocate cache over a write-through/no-write-allocate cache?

4. In Intel processors, caches use physical addresses, whereas instructions use virtual addresses. How is it then possible that L1 cache is still efficiently accessed?

5. Consider a computer with a direct-mapped TLB cache with 8 total entries. The page size is 1024 bytes. In the following code, assume that the vector  $x$  starts at address 0 in memory, and that  $n = 128$ . Determine memory locations (addresses) of  $y[0]$  and  $z[0]$  that would lead to the worst TLB hit-rate. State also this worst TLB hit-rate and justify your answer. Assume that expressions are evaluated from left to right.

```
1 void compute(float* x, float* y, float* z, int n){
2     for(int i=0; i < n; i++){
3         x[i] = x[i] + y[i] + z[i];
4     }
5 }
```

